

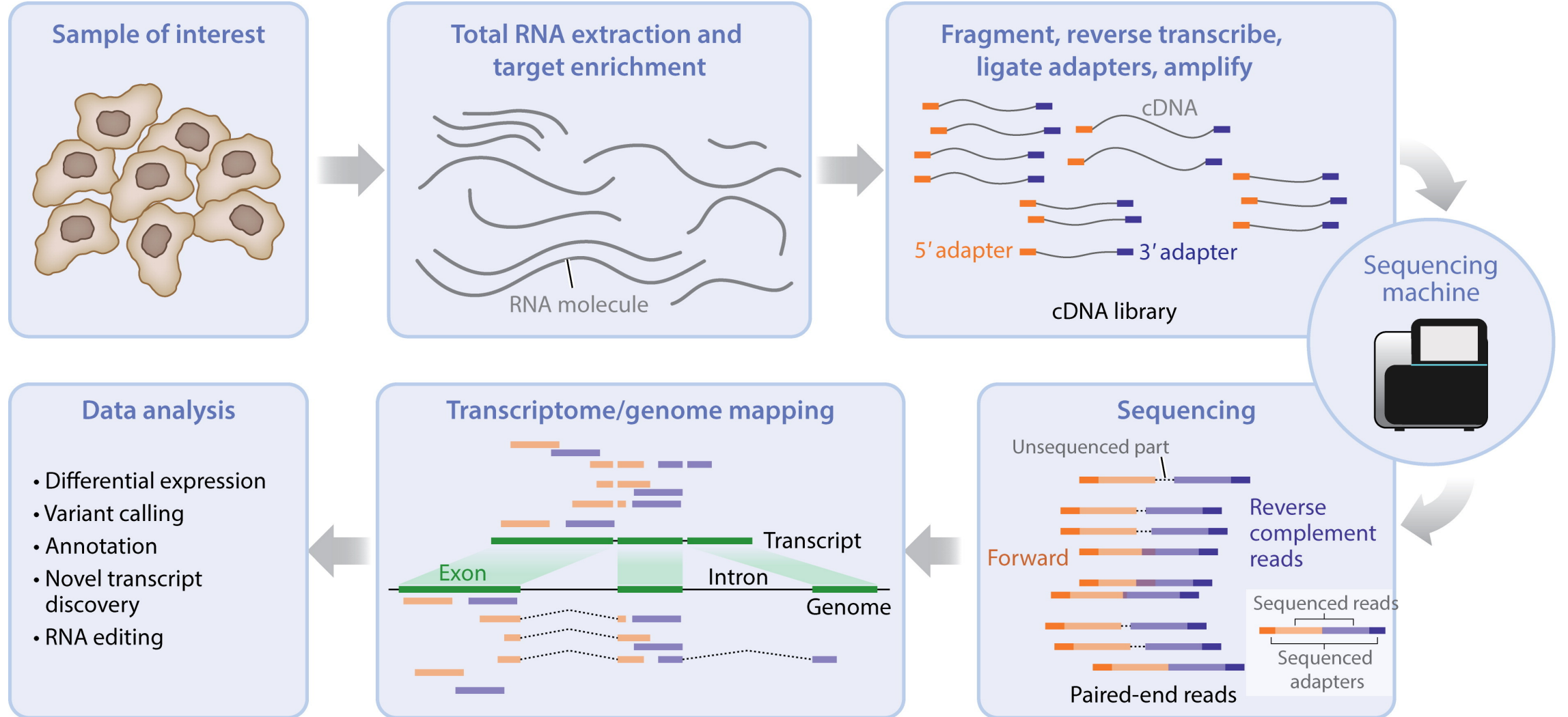
Tools For Computational Biology – Lecture 17

Introduction to (bulk) RNA-seq

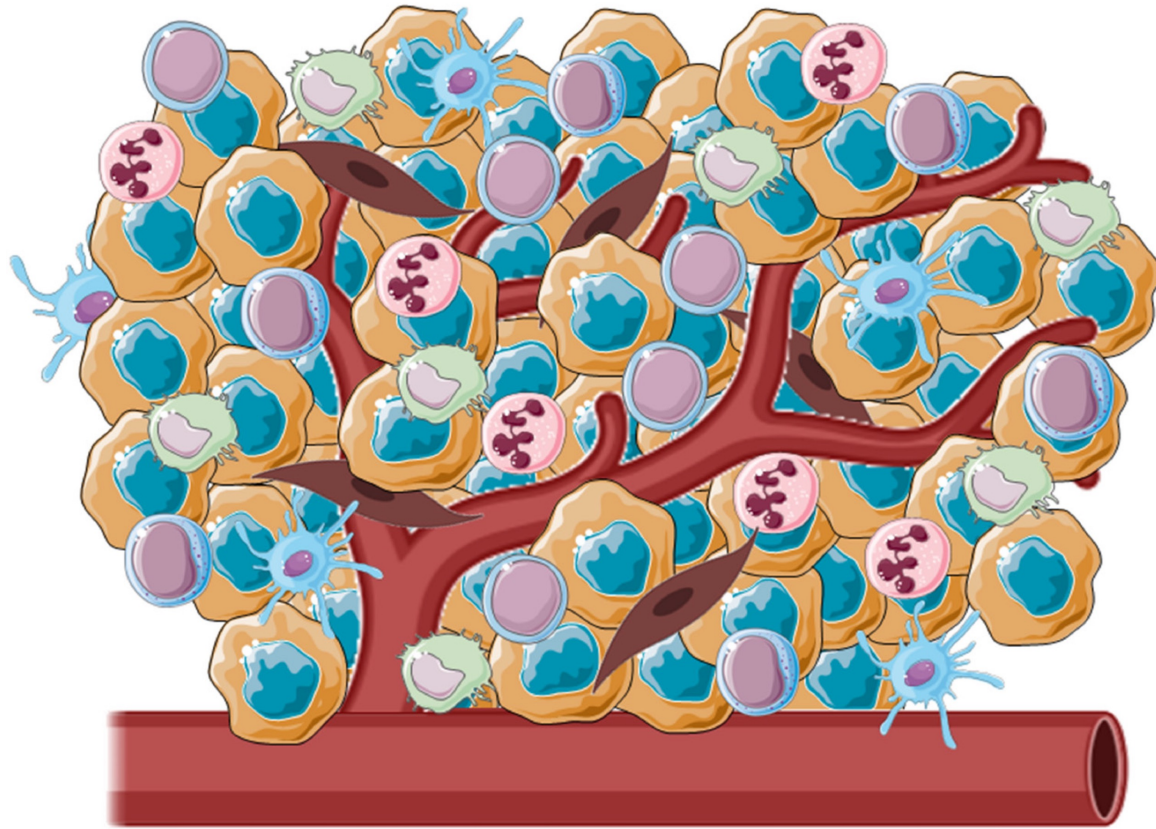
Overview

- A very brief introduction to RNA-seq
- Bulk RNA-seq analysis – R markdown
 - Quality controls and checks
 - Count matrices
 - Differential analysis
 - Gene Set Enrichment Analysis / Gene Ontology Analysis
 - Visualization

RNA-seq



RNA-seq



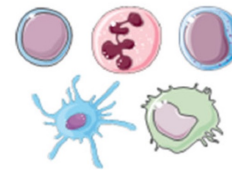
Tumor cells



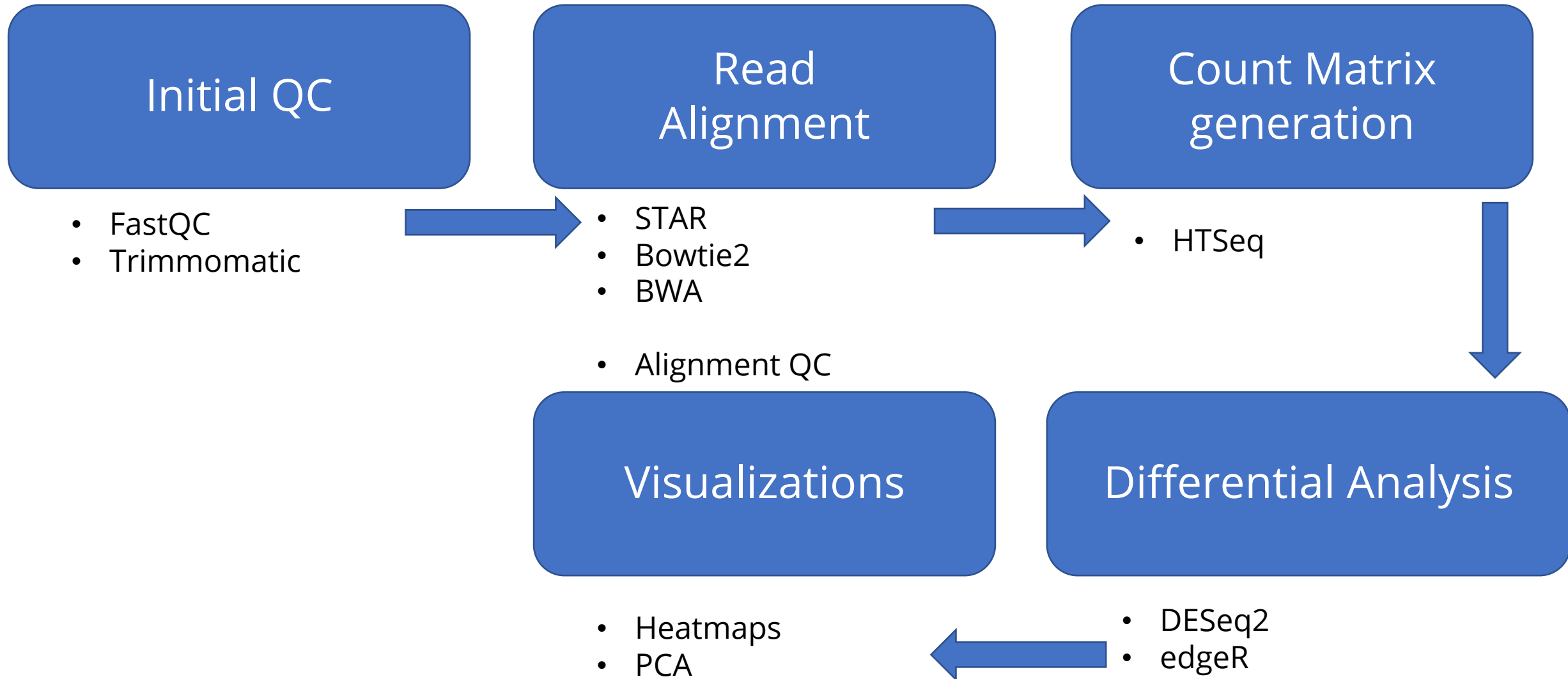
Cancer-associated fibroblast



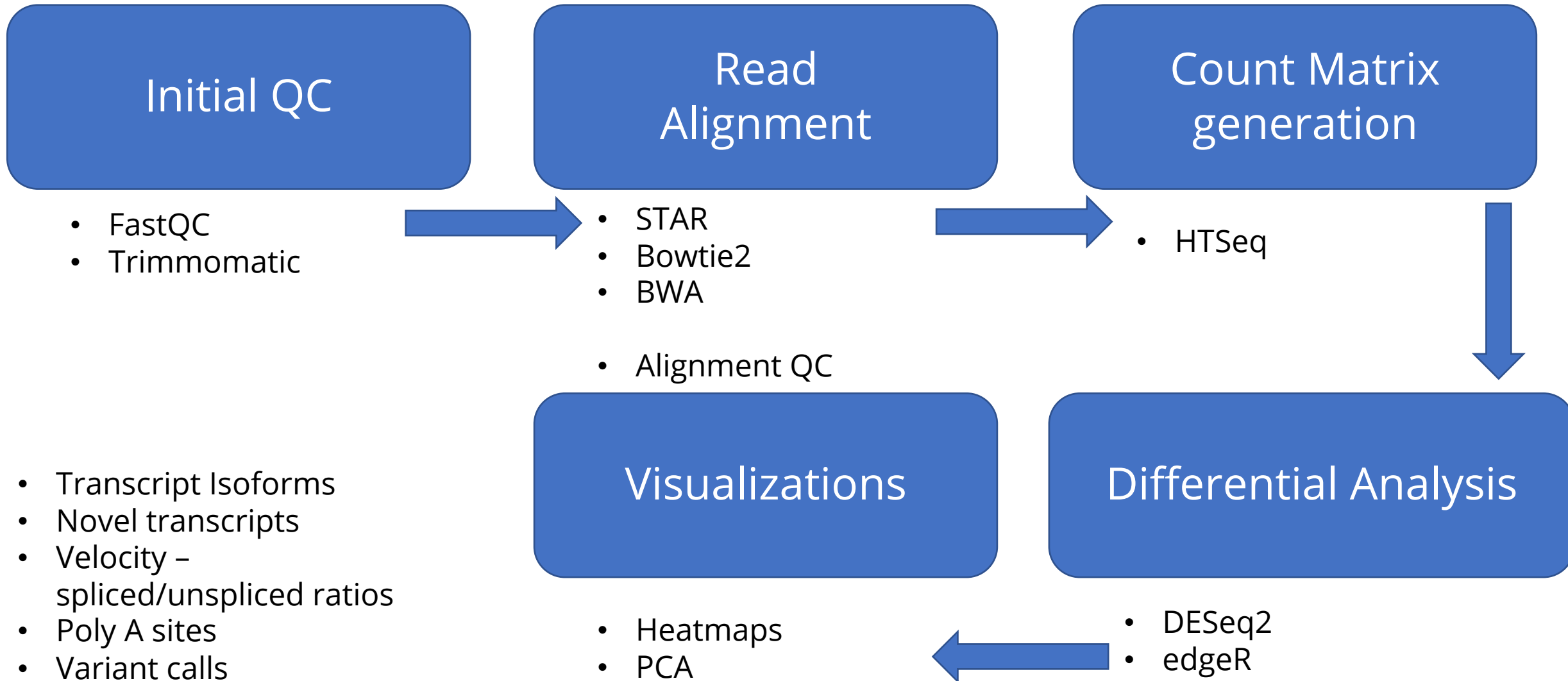
Immune cells



RNA-seq analysis steps



RNA-seq analysis steps



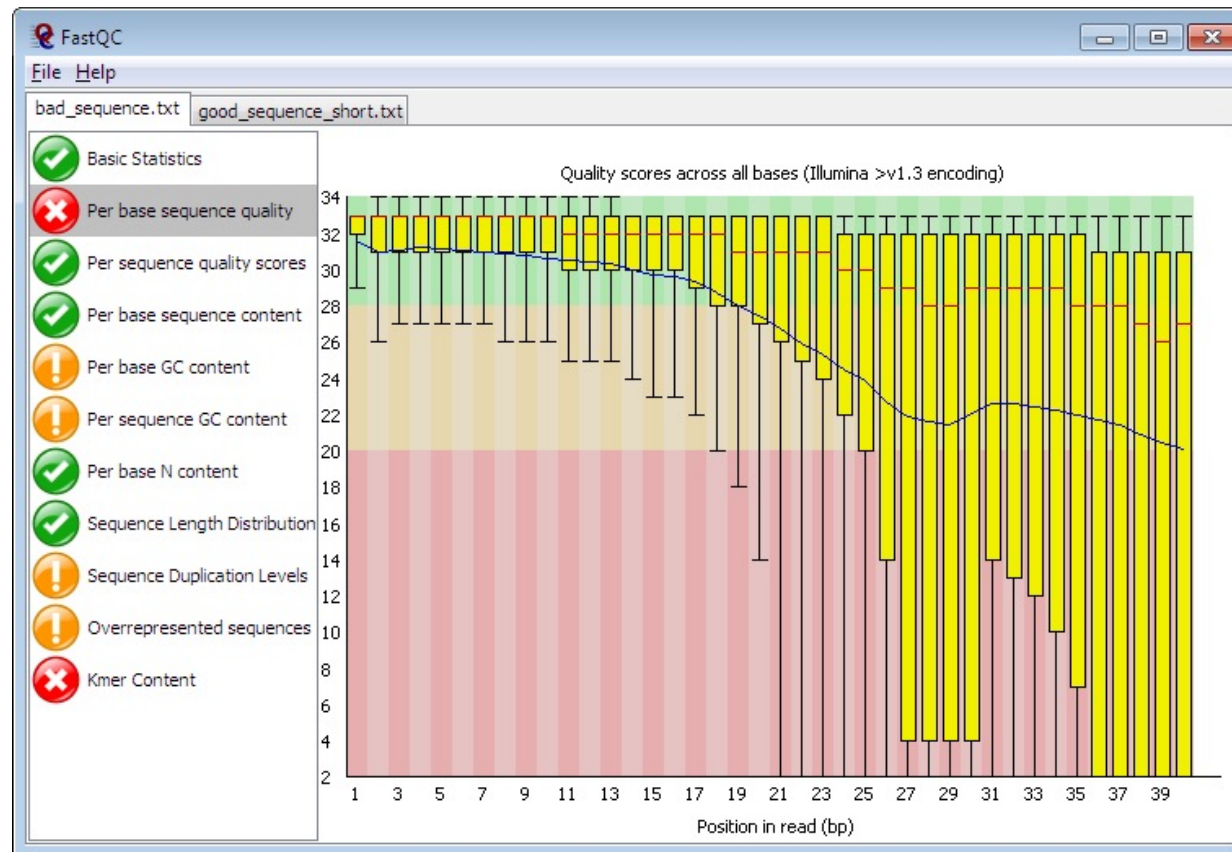
Initial QC: Fastq files

- Fastq File format

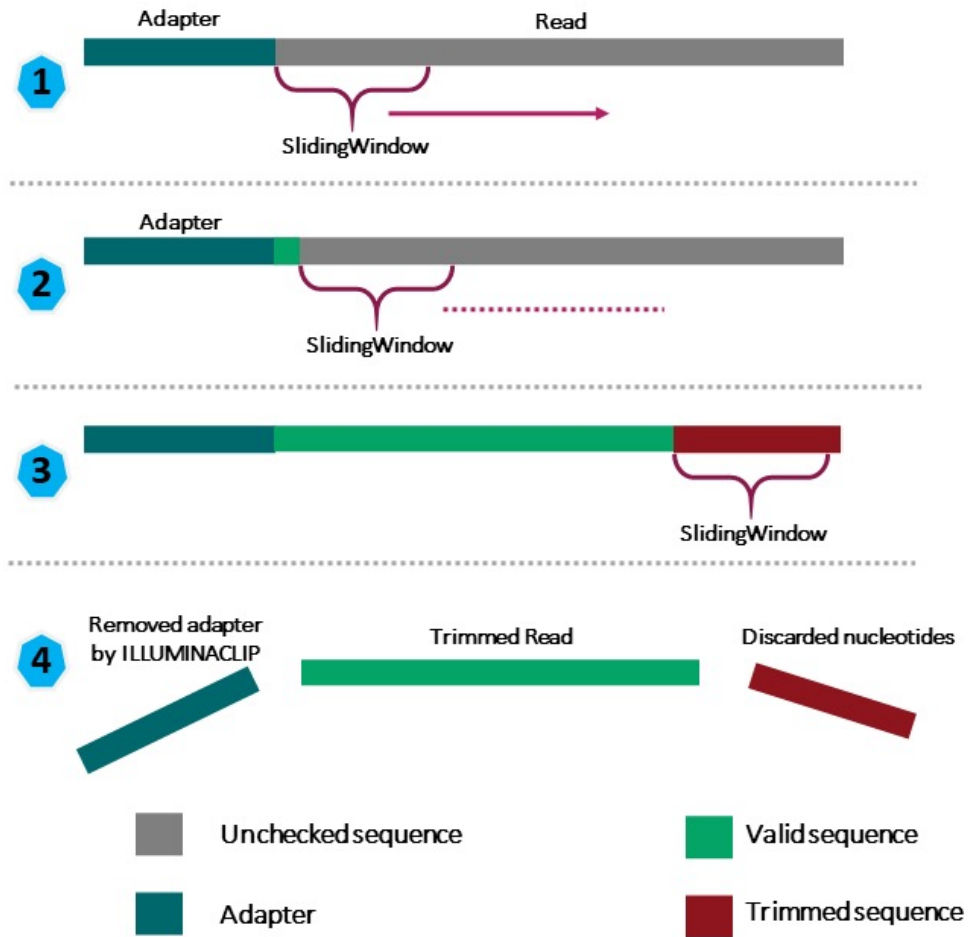
```
Read N @K00217:70:HHJF5BBXX:1:1101:1955:1121 2:N:0:NTGCGCGC  
NAAAAGAAAAGAACCCGCCGAGCAGTCAAATTCCAGAGGGCCATTACTGG  
+  
#A<---FJJJFJJJJ-AAJFJFA-AFJ<-FF<F7J<7--7A-7--<7-7-  
  
Read N + 1 @K00217:70:HHJF5BBXX:1:1101:2016:1121 2:N:0:NGCGAGTA  
NTCTGTCACGCACATGTGTCCTGTGGGTATAGCTAGAAGGACAGGAGGCT  
+  
#-<<7FJAJAFJJJJ7AJA<F-FJ<J-7F-FJJ-<AJJFFJFF-7-7--7
```

FastQC

- Quality check for Fastq files
- <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>



Trimmomatic: Remove sequencing adapters



Alignment: STAR

- STAR Aligner: Fast!

Aligner	Mapping speed: million read pairs/hour		Peak physical RAM, GB	
STAR	309.2	549.9	27.0	28.4
STAR sparse	227.6	423.1	15.6	16.0
TopHat2	8.0	10.1	4.1	11.3
RUM	5.1	7.6	26.9	53.8
MapSplice	3.0	3.1	3.3	3.3
GSNAP	1.8	2.8	25.9	27.0

- Map to both genome & transcriptome
- Output: Bam files with information about genomic location of reads

Alignment QC

- STAR log

```
Started job on | Apr 23 23:17:02
Started mapping on | Apr 23 23:17:04
Finished on | Apr 23 23:26:52
Mapping speed, Million of reads per hour | 115.68
```

```
Number of input reads | 18894432
Average input read length | 298
```

UNIQUE READS:

```
Uniquely mapped reads number | 17704240
Uniquely mapped reads % | 93.70%
```

```
Average mapped length | 297.39
Number of splices: Total | 3119841
Number of splices: Annotated (sjdb) | 2663436
Number of splices: GT/AG | 3080422
Number of splices: GC/AG | 14219
Number of splices: AT/AC | 248
Number of splices: Non-canonical | 24952
Mismatch rate per base, % | 0.49%
Deletion rate per base | 0.02%
Deletion average length | 2.70
Insertion rate per base | 0.02%
Insertion average length | 2.30
```

MULTI-MAPPING READS:

```
Number of reads mapped to multiple loci | 806405
% of reads mapped to multiple loci | 4.27%
Number of reads mapped to too many loci | 1146
% of reads mapped to too many loci | 0.01%
```

Count matrix generation

- Count matrix: Gene expression counts for each sample
- We need two components
 - BAM files
 - ???

Count matrix generation

- Count matrix: Gene expression counts for each sample
- We need two components
 - BAM files
 - Transcript definition!

